

1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025 2026 2027 2028 2029 2030 2031 2032 2033 2034 2035 2036 2037 2038 2039 2040 2041 2042 2043 2044 2045 2046 2047 2048 2049 2050 2051 2052 2053 2054 2055 2056 2057 2058 2059 2060 2061 2062 2063 2064 2065 2066 2067 2068 2069 2070 2071 2072 2073 2074 2075 2076 2077 2078 2079 2080 2081 2082 2083 2084 2085 2086 2087 2088 2089 2090 2091 2092 2093 2094 2095 2096 2097 2098 2099 2100 2101 2102 2103 2104 2105 2106 2107 2108 2109 2110 2111 2112 2113 2114 2115 2116 2117 2118 2119 2120 2121 2122 2123 2124 2125 2126 2127 2128 2129 2130 2131 2132 2133 2134 2135 2136 2137 2138 2139 2140 2141 2142 2143 2144 2145 2146 2147 2148 2149 2150 2151 2152 2153 2154 2155 2156 2157 2158 2159 2160 2161 2162 2163 2164 2165 2166 2167 2168 2169 2170 2171 2172 2173 2174 2175 2176 2177 2178 2179 2180 2181 2182 2183 2184 2185 2186 2187 2188 2189 2190 2191 2192 2193 2194 2195 2196 2197 2198 2199 2200 2201 2202 2203 2204 2205 2206 2207 2208 2209 2210 2211 2212 2213 2214 2215 2216 2217 2218 2219 2220 2221 2222 2223 2224 2225 2226 2227 2228 2229 2230 2231 2232 2233 2234 2235 2236 2237 2238 2239 2240 2241 2242 2243 2244 2245 2246 2247 2248 2249 2250 2251 2252 2253 2254 2255 2256 2257 2258 2259 2260 2261 2262 2263 2264 2265 2266 2267 2268 2269 2270 2271 2272 2273 2274 2275 2276 2277 2278 2279 2280 2281 2282 2283 2284 2285 2286 2287 2288 2289 2290 2291 2292 2293 2294 2295 2296 2297 2298 2299 2300 2301 2302 2303 2304 2305 2306 2307 2308 2309 2310 2311 2312 2313 2314 2315 2316 2317 2318 2319 2320 2321 2322 2323 2324 2325 2326 2327 2328 2329 2330 2331 2332 2333 2334 2335 2336 2337 2338 2339 2340 2341 2342 2343 2344 2345 2346 2347 2348 2349 2350 2351 2352 2353 2354 2355 2356 2357 2358 2359 2360 2361 2362 2363 2364 2365 2366 2367 2368 2369 2370 2371 2372 2373 2374 2375 2376 2377 2378 2379 2380 2381 2382 2383 2384 2385 2386 2387 2388 2389 2390 2391 2392 2393 2394 2395 2396 2397 2398 2399 2400 2401 2402 2403 2404 2405 2406 2407 2408 2409 2410 2411 2412 2413 2414 2415 2416 2417 2418 2419 2420 2421 2422 2423 2424 2425 2426 2427 2428 2429 2430 2431 2432 2433 2434 2435 2436 2437 2438 2439 2440 2441 2442 2443 2444 2445 2446 2447 2448 2449 2450 2451 2452 2453 2454 2455 2456 2457 2458 2459 2460 2461 2462 2463 2464 2465 2466 2467 2468 2469 2470 2471 2472 2473 2474 2475 2476 2477 2478 2479 2480 2481 2482 2483 2484 2485 2486 2487 2488 2489 2490 2491 2492 2493 2494 2495 2496 2497 2498 2499 2500 2501 2502 2503 2504 2505 2506 2507 2508 2509 2510 2511 2512 2513 2514 2515 2516 2517 2518 2519 2520 2521 2522 2523 2524 2525 2526 2527 2528 2529 2530 2531 2532 2533 2534 2535 2536 2537 2538 2539 2540 2541 2542 2543 2544 2545 2546 2547 2548 2549 2550 2551 2552 2553 2554 2555 2556 2557 2558 2559 2560 2561 2562 2563 2564 2565 2566 2567 2568 2569 2570 2571 2572 2573 2574 2575 2576 2577 2578 2579 2580 2581 2582 2583 2584 2585 2586 2587 2588 2589 2590 2591 2592 2593 2594 2595 2596 2597 2598 2599 2600 2601 2602 2603 2604 2605 2606 2607 2608 2609 2610 2611 2612 2613 2614 2615 2616 2617 2618 2619 2620 2621 2622 2623 2624 2625 2626 2627 2628 2629 2630 2631 2632 2633 2634 2635 2636 2637 2638 2639 2640 2641 2642 2643 2644 2645 2646 2647 2648 2649 2650 2651 2652 2653 2654 2655 2656 2657 2658 2659 2660 2661 2662 2663 2664 2665 2666 2667 2668 2669 2670 2671 2672 2673 2674 2675 2676 2677 2678 2679 2680 2681 2682 2683 2684 2685 2686 2687 2688 2689 2690 2691 2692 2693 2694 2695 2696 2697 2698 2699 2700 2701 2702 2703 2704 2705 2706 2707 2708 2709 2710 2711 2712 2713 2714 2715 2716 2717 2718 2719 2720 2721 2722 2723 2724 2725 2726 2727 2728 2729 2730 2731 2732 2733 2734 2735 2736 2737 2738 2739 2740 2741 2742 2743 2744 2745 2746 2747 2748 2749 2750 2751 2752 2753 2754 2755 2756 2757 2758 2759 2760 2761 2762 2763 2764 2765 2766 2767 2768 2769 2770 2771 2772 2773 2774 2775 2776 2777 2778 2779 2780 2781 2782 2783 2784 2785 2786 2787 2788 2789 2790 2791 2792 2793 2794 2795 2796 2797 2798 2799 2800 2801 2802 2803 2804 2805 2806 2807 2808

BLAKELY, SOKOLOFF, TAYLOR & ZAFMAN LLP
12400 Wilshire Blvd., 7th Floor
Los Angeles, CA 90025-1026
(714) 557-3800

A PLATFORM AND METHOD FOR REPRESENTING AND
SUPPORTING HOT-PLUGGED NODES

1. Field

This invention relates to the field of computers. In particular, the invention relates to a platform employing a mechanism for representing and supporting a hot-plugged node and its constituent components as a collective unit to its operating system.

2. Background

Advances in technology have opened up many opportunities for applications that go beyond the traditional ways of doing business. Electronic commerce over the Internet has become widely accepted, requiring many companies to either install one or more servers to host a web site and maintain accessible databases or to contract with data centers to provide such services. In addition to performance, important functional characteristics for these servers include reliability, availability and serviceability.

Normally, conventional server architectures feature both processors and memory coupled to a front-side bus. This conventional server architecture greatly hinders server performance due to a number of factors. For instance, one factor is that the front-side bus is non-scalable. Thus, as more processors and memory have access to the front-side bus, bandwidth constraints associated with that bus adversely effect overall server performance. Multi-node architectures where processors, memory and input/output (I/O) components are distributed across multiple interconnected nodes overcomes the limitations of front-side bus and enables building larger systems with scalable server performance.

Another factor is that any node, namely a primary component interconnected to a group of components (referred to as "constituent components"), may be hot-plugged

to allow its addition or removal while the operating system (OS) of the server continues to operate. In order to provide a hot-plug solution, however, the constituent components must be represented and visible to the OS of the server. One option is to represent the constituent components, such as one or more processors and memories for
5 example, separately in accordance with a platform firmware interface such as the Advanced Configuration and Power Interface (ACPI) Specification (Version 2.0) published July 27, 2000. However, for those platforms supporting non-uniform memory access (NUMA) architectures, this separate representation poses a number of disadvantages.

10 For instance, one disadvantage is that the OS would not be able to determine proximity relationships between nodes. In other words, the OS would not be able to determine which processor(s) and which memory are interconnected to the same node and adjacent to each other. Such proximity data, if available, would allow the OS to attempt to allocate memory for a processor from the same node in order to avoid time
15 latency penalties caused by accessing memory that is remotely located from that processor.

Another disadvantage is that, during removal of a hot-plugged node, the OS is unable to simultaneously determine which components were removed with the node. Also, during addition of a hot-plugged node inclusive of a processor and memory, the
20 OS of the server may initialize the processor prior to activation of the memory. Hence, the processor specific memory allocation may be inefficient because remotely located memory would be allocated to the processor before local memory is available. This would adversely impact overall server performance and further complicate removal of hot-plugged nodes.

25 Also, it is contemplated that processor(s) and memory of a node must be initialized to a known state before the OS is made aware of them. When the node is

hot-plugged and separate from the substrate maintaining the Basic Input/Output System (BIOS), the BIOS cannot be used for initialization of the processor(s) and memory.

One reason is that when the platform under control of the OS, the OS only allows nodes to initiate non-coherent transactions until recognized by the OS. However, the hot-

- 5 plugged node need to initiate coherent (memory) transactions to reprogram registers to
enable various communication ports and links.

BRIEF DESCRIPTION OF THE DRAWINGS

The features and advantages of the present invention will become apparent from the following detailed description of the present invention in which:

Figure 1 is a first exemplary embodiment of a substrate layout for a platform
5 utilizing the invention.

Figure 2 is a second exemplary embodiment of a substrate layout of a platform utilizing the invention.

Figure 3 is an exemplary embodiment of a scalability node controller implemented within the platforms of Figures 1-2.

10 Figure 4 is an exemplary embodiment of a scalability port switch implemented within the platforms of Figures 1-2.

Figure 5 is an exemplary embodiment of the platform of Figure 2 prior to undergoing dynamic partitioning.

15 Figure 6 is an exemplary embodiment of the platform of Figures 2 after undergoing dynamic partitioning to produce two 4-way platforms.

Figure 7 is an exemplary embodiment of representation of a container object that is representative of a node and generally consistent with ACPI standards.

DESCRIPTION

The invention relates to a platform and method for supporting and representing a hot-plugged node and its constituent components as a collective unit to its operating system (OS). The constituent components are initialized to a known state through a distributed BIOS mechanism before the OS is made aware of their presence. For this embodiment, the platform is configured with OS-directed device configuration and power management of both the node/components and the platform itself.

Herein, certain details are set forth in order to provide a thorough understanding of the invention. It is apparent to a person of ordinary skill in the art, however, that the invention may be practiced through many embodiments other than those illustrated. Well-known circuits and ACPI parameters are not set forth in detail in order to avoid unnecessarily obscuring the invention.

In the following description, terminology is used to discuss certain features of the present invention. For example, a "platform" includes hardware equipment and/or software that process data. One type of platform is a computer such as a server, although other types of hardware equipment may employ aspects of the invention. A "code segment" is a piece of software (e.g., code or instructions) that, when executed, performs a certain function. Software code is stored in platform readable medium, which may include any medium that can store or transfer information. Examples of the platform readable medium include an electronic circuit, a semiconductor memory device, a volatile memory (e.g., random access memory "RAM"), a non-volatile memory (e.g., read-only memory, a flash memory, etc.), a floppy diskette, a compact disk, an optical disk, a hard drive disk, a fiber optic medium, and the like.

In addition, a "device" is a single component or a collection of interconnected components that is also referred to as a "node". Herein, a node may be referenced and identified by its primary component. Each component may be an active device (e.g.,

integrated circuit, timing or clocking components, etc.), but it is contemplated that the invention may be applicable for use with passive components (e.g., resistors, capacitors, inductors, etc.). A “link” is broadly defined as any type of information-carrying medium such as electrical wire, optical fiber, cable, trace bus or even wireless signaling technology. In addition, the term “hot-plug” or any tense thereof indicates a characteristic where a device may be added, removed or replaced while the OS of the platform continues to operate.

I. Platform Hardware Architecture Overview with Distributed BIOS

Referring to Figure 1, a first exemplary embodiment of a platform utilizing the invention is shown. The platform 100 comprises a processor substrate 110, an input/output (I/O) substrate 170 and an interconnection substrate 150 that couples devices mounted on the processor substrate 110 with those on the I/O substrate 170. Each “substrate” is formed from any type of material or combination of materials upon which integrated circuits as well as a wide variety of other types of devices (e.g., passive, sockets, timing, etc.) can be attached. Each substrate may be produced in a number of form factors such as, for example, a circuit board acting as a motherboard or a removable daughter card.

As shown, the processor substrate 110 comprises a first scalable node controller (SNC0) 120 that is configured with hot-plug capability as shown in Figure 2. SNC0 120 is coupled to a connector 115 placed at an edge of the substrate 110. This connector 115 is adapted for coupling with a mating connector 155 placed on the interconnection substrate 150. SNC0 120 is further coupled to a processor cluster 125 supporting one or more processors 127₁-127_M (“M” being a positive integer), a local memory cluster 130 having one or more banks of memory 133 and a firmware hub 140. The firmware hub 140 is configured to store Basic Input/Output System (BIOS) 141 configured for partial initialization of components and enablement of links therefrom as described in Figure 2 (hereinafter referred to as “INIT BIOS”).

Referring still to Figure 1, SNC0 120 features two scalability port interfaces 124 (see also Figure 3) that are both coupled to connector 115 via links 160 and 165. This enables data to be routed from SNC0 120 to a Server Input/Output Hub (SIOH) 180 via connectors 115 and 155 as well as connector 175 of I/O substrate 170. SIOH0 180 provides communications with high-speed links. For example, SIOH0 180 provides coupling to one or more bridges 185 (e.g., P64H2 devices) that support communications with one or more I/O buses such as a Peripheral Component Interconnect “PCI” bus and/or a higher speed PCI bus which is referred to as the “PCI-X bus” for example. SIOH0 180 further provides coupling to a virtual interface bridge (VXB) 190 (also referred to as “host channel adapter”) and an I/O Riser substrate having an input/output control hub (ICH2) 196 mounted thereon. The VXB 190 provides a four 10-bit system I/O full-duplex channels. ICH2 196 supports a number of functions that are designed to support platform security in addition to traditional I/O and platform boot functions. ICH2 196 enables communications with a boot flash containing a system BIOS for booting the platform (not shown), networking ports as well as various I/O peripherals such as a mouse, alphanumeric keyboard, and the like (not shown).

Referring now to Figure 2, a second exemplary embodiment of a multi-node platform utilizing the invention is shown. Platform 200 is configured to support multiple processor substrates that enable the M-way processor-based platform 100 of Figure 1 to be converted to the M+N-way platform 200 as shown. For this embodiment, as shown, platform 200 comprises first processor substrate 110 and a second processor substrate 210, both coupled to a multi-substrate interconnection substrate 250. The dual-substrate interconnection substrate 250 is coupled to an I/O substrate 270.

More specifically, as shown in both Figures 2 and 3, first processor substrate 110 comprises SNC0 120 coupled to processor cluster 125, local memory cluster 130,

firmware hub 140 and connector. SNC0 120 comprises a plurality of port interface that, when activated, enable communications over different links. For example, a processor port interface 121 of SNC0 120 provides a communication path to processors 127₁-127_M of processor cluster 125 via processor link 126. Memory port interface 122 of SNC0 120 provides a communication path to local memory cluster 130 via a memory link 131. In one embodiment, memory link 131 provides four communication sub-links 132₁-132₄ supporting a total data throughput of approximately 6.4 Gigabytes per second (GB/s). Each of the sub-links 132₁,..., 132₄ may be coupled to a bank of local memory devices 133 (e.g., RDRAM) or a memory repeater hub 134₁,..., 134₄ that operates as an RDRAM-to-SDRAM translation bridge.

SNC0 120 further includes a first scalability port interface 124₁ that enables a communication path over link 260 to a first scalability port switch (SPS0) 275 via connector 115 and mating connector 255. SNC0 also includes a second scalability port interface 124₂ that enables a communication path over link 261 to a second scalability port switch (SPS1) 276 via connectors 115 and 255.

As further shown in Figures 2 and 3, SNC0 120 comprises a port interface 123 that enables a communication path to firmware hub 140 via link 142. Firmware hub 140 comprises INIT BIOS 141 that is configured to initialize processors 127₁-127_M, local memory 133, and scalability port interfaces 124₁ and 124₂ to communicate with the OS. As a result, the distributed INIT BIOS 141 enables hot-plug addition of a boot node, namely first processor substrate 110, and supports dynamic partitioning of platform 200.

Similar in architecture to first processor substrate 110, second processor substrate 210 comprises a second scalable node controller (SNC1) 220 that is mounted on a substrate and coupled to a processor cluster 211, a local memory cluster 216, a firmware hub 240 as well as a connector 215. Connector 215 is adapted to couple with a second mating connector 256 of interconnection substrate 250.

As shown in Figure 2, processor cluster 211 comprises a processor link 212 interconnecting one or more processors 213₁-213_N ("N" being a positive integer). It is contemplated that these N processors may equal in number to the M processors provided by first processor substrate 110, although such a 1:1 correlation is not necessary. Processor cluster 211 is coupled to a processor port interface of SNC1 220 via processor link 212. Local memory cluster 216 is coupled to a memory port interface of SNC1 220 through a memory link 217. SNC1 220 features two scalability port interfaces 221 that are both coupled to connector 215 via links 222 and 223.

As further shown in Figure 2, SNC1 220 comprises a port interface 224 that enables a communication path to firmware hub 240 via link 242. Firmware hub 240 comprises INIT BIOS 241 that is configured to initialize processors 213₁-213_N, local memory 218, and scalability port interfaces 222 and 223 to support communications with the OS when a hot-plugged operation occurs involving second processor substrate 110. The portion of INIT BIOS 241 enables hot-plug addition of another boot node (e.g., second processor substrate 210) and also supports dynamic partitioning of platform 200 as described in connection with Figures 5 and 6.

Referring still to Figure 2, interconnection substrate 250 enables data to be propagated from SNC0 120 to both SPS0 275 and SPS1 276. In particular, first mating connector 255 receives data transferred through connector 115 and propagates that data over links 260 and 261. Links 260 and 261 are coupled to a connector 265 of interconnection substrate 250. Connector 265 may be coupled to a mating connector 271 of I/O substrate 270, which propagates the data from links 260 and 261 to SPS0 275 and SPS1 276, respectively. Similarly, in a redundant fashion, interconnection substrate 250 enables data to be propagated from SNC1 220 to SPS0 275 and SPS1 276 over links 262 and 263, respectively.

As shown in Figure 4, in one embodiment, SPS0 275 and/or SPS1 276 is a crossbar switch (e.g., integrated 6x6 crossbar) that enables communication with

components over six port interfaces 300-305. For example, with this embodiment, each scalability port switch would enable communications between four SNCs and two SIOHs. Both SPS0 275 and SPS1 276 are programmed by accessing internal control and status registers via PCI configuration interface, System Management Bus (SMBus) interface, or Joint Test Action Group (JTAG) interface.

Referring back to Figure 2, I/O substrate 270 comprises SPS0 275 and SPS1 276, each coupled to a first Server Input/Output Hub (SIOH0) 280 and a second Server Input/Output Hub (SIOH1) 285. As previously described, both SIOH0 280 and SIOH1 285 provide communications with high-speed links. For example, SIOH1 285 provides coupling to one of more of the following: (1) one or more bridges 290 (e.g., P64H2 devices) that support communications with one or more I/O buses; (2) a virtual interface bridge (VXB) 291 that provides system I/O full-duplex channels; and/or (3) an I/O Riser substrate 292 having an input/output control hub (ICH2) 293 mounted thereon.

Besides system BIOS software retrieved via ICH2, various portions of INIT BIOS 141 are configured to reside in firmware hub 140 being coupled to SNC0 120. Likewise, various portions of the INIT BIOS 241 reside in firmware hub 240. As shown herein, both INIT BIOS 141 and/or 241 are implemented in a distributed fashion to assist in initialization without leaving such responsibility to the OS.

Herein, both INIT BIOS 141 and 241 may be responsible for electing their respective node boot strap processor to handle initial BIOS boot sequence operations for its specific substrate and enable communication path(s) to SPS0 275 and/or SPS1 276. For instance, INIT BIOS 141 enables the scalability port interfaces 124 and waits for Idle Flits from SPS0 275 and/or SPS1 276. Likewise, INIT BIOS 241 enables the scalability port interfaces 221 and waits for Idle Flits from SPS0 275 and/or SPS1 276.

During a normal boot of the platform 200, the node boot strap processors in turn elect the system boot strap processor which runs the system BIOS located in its respective FWH attached to ICH2 in order to complete the boot of the platform. Both INIT BIOS 141 and/or 241 is further configured to initialize the processors and memory on a hot-plugged node. For instance, INIT BIOS 141 is configured to initialize processors and memory associated with SNC0 120, which requires SNC0 120 to read the configuration state information from SPS0 275 and SPS1 276 using non-coherent accesses. Additionally, both INIT BIOS 141 and 241 program registers to indicate resources of the hot-plugged node (e.g., memory interleave registers to indicate memory of the new node) and to notify the OS of the presence of a fully initialized hot-plugged node.

For an 8-way platform featuring processor substrates 110 and 210 as shown, after successful boot of platform 200, the user may decide to remove various resources. For instance, platform 200 may undergo a hot-plug removal of a node featuring SCN0 120 mounted on first processor substrate 110. This may be accomplished by the OS transmitting an ejection notice for a container object that identifies SCN0 as well as its constituent components coupled thereto (see Figure 7). Likewise, when undergoing a hot-plug addition of a node (e.g., SCN0 120 being the primary component), after initialization of its constituent components by the distributed INIT BIOS 141, the OS would bring local memory 133 online prior to processors 127₁-127_M so that such memory may be allocated to processors 127₁-127_M before selecting remotely located memory.

Dynamic partitioning can be defined as an ability to either split one R-way platform ("R" being a positive integer) into multiple smaller systems while the original OS continues to run without shutdown or merge multiple partitions into one larger partition while the original OS continues to run without shutdown. For example, using 8-way platform 200 of Figure 2 for illustrative purposes, dynamic partitioning allows 8-

way platform 200 to be split into two 4-way platforms 400 and 500 or an ability to merge two 4-way platforms 400 and 500 to form single 8-way platform 200. The dynamic partitioning operation occurs without requiring the OS to shutdown, namely reboot.

Dynamic partitioning is accomplished by using the hot-plug capability. As shown in Figures 5 and 6, for this embodiment, in order to split 8-way platform 200 into 4-way platforms 400 and 500, the following sequence of operations are performed:

1. Indicate a hot-plug removal event to the OS for SIOH1 285;
2. Indicate a hot-plug removal event to the OS for CPU/Memory node1 (e.g., second processor substrate 210);
3. When both hot-plug removal operations are complete, the original OS is running on 4-way platform 400 comprising CPU/Memory node0 (first processor substrate 110) and SIOH0 280;
4. Program registers associated with SPS0 275 and SPS1 276 to indicate that the platform is partitioned into two;
5. Initialize the new platform 500 comprising CPU/Memory node1 (second processor substrate 210) and SIOH1 285. This platform 500 is able to run its own copy of OS and applications independent of platform 400.

Similarly, for merging of two 4-way platforms 400 and 500 into 8-way platform 200, the new nodes are announced to the running OS as hot-plug events. As a result, OS is able to add the hot-plugged CPU/Memory node and an I/O node (e.g., SIOH and constituent components) to the running platform 200 without any interruption to its service.

One usage of a dynamic partitioning is for rolling upgrades. Applications, other software services and even the OS require updates from time to time to either enhance functionality or to fix existing problem. Typically, in high reliability mission critical environments, the software updates should not be applied directly to the targeted platform in the field. However, testing the software update on a different server may not accurately test the real world environment in which it is deployed. Hence, the running platform is split into two using dynamic domain partitioning and the software update is applied to the newly formed partition while the original OS and software continues to run on the other partition. After sufficient testing, the partition running the old software is merged into the partition running the updated software thus accomplishing the update in an efficient and reliable manner.

II. Container Object Representation

Referring not to Figure 7, an exemplary embodiment of representation of a container object that is representative of a node (e.g., SNC0 being interconnected to constituent components such as processors and memory) and generally consistent with ACPI standards is shown. Herein, container object 600 provides a mechanism for handling an ejection notice for hot-plug removal of SNC0 from the platform. This container object 600 provides at least information as to those devices that are constituent components of SNC0.

Herein, container object 600 of the first scalability node controller identifies its constituent components. For example, container object 600 represents SNC0 as a node including processors and memory devices as constituent components of that node. The eject (_EJ0) method 610 is invoked to eject, during a hot-plug removal of the first processor substrate or perhaps SNC0, those constituent components of SNC0. This occurs before SNC0 is ejected. For this embodiment, since SNC0 is coupled to processors and memory, container object 600 comprises a hardware identification (_HID) object 620, a proximity (_PXM) object 630, a processor object 640, and a

device object 650. Each processor or device object 640 and 650 may be generally referred to as a "component object."

As shown in Figure 7, the _HID object 620 contains a string for identifying the device type associated with the container object to the OS for power management and configuration. As an example, for this embodiment, the _HID object 610 would return the value of "ACPI0004" to identify SCN0 as a node.

The _PXM object 630 is used to describe proximity domains (i.e., groupings of devices) within the platform. In other words, the _PXM object 630 provides an integer that identifies a device as belonging to a specific proximity domain. The OS assumes that two devices in the same proximity domain are coupled and tends to allocate memory to those processors within the same proximity domain as the memory. For instance, SNC0 may be assigned "0" to denote that it is in a first domain while SNC1 would be assigned "1" to denote that it is in a secondary domain.

The processor object 640 is used to identify which processor(s) constitute components associated with SNC0. Similarly, device object 650 is used to identify other devices (e.g., memory) that constitute components associated with SNC0.

By using the container object representation, the OS can easily determine the components that belong together and thereafter is able to optimize memory allocations and other operations. Moreover, in the case of a hot-plug removal, the ejection notice is just sent to the node, which in turn propagates down to the constituent components. For example, in case of a SNC, the alternative would be to notify each processor and memory component individually each of which would be interpreted by the OS as an individual hot-plug removal operation thereby making the hot-plug removal operation inefficient. Also, during a hot-plug addition of a node, the OS brings the associated memory online before the processors are brought in attempts to allocate memory that is local to the processors.

